# FIRM: **F**air and H**I**gh-Perfo**R**mance **M**emory Control for Persistent Memory Systems

Jishen Zhao     Onur Mutlu     Yuan Xie

# New Design Opportunity with NVMs
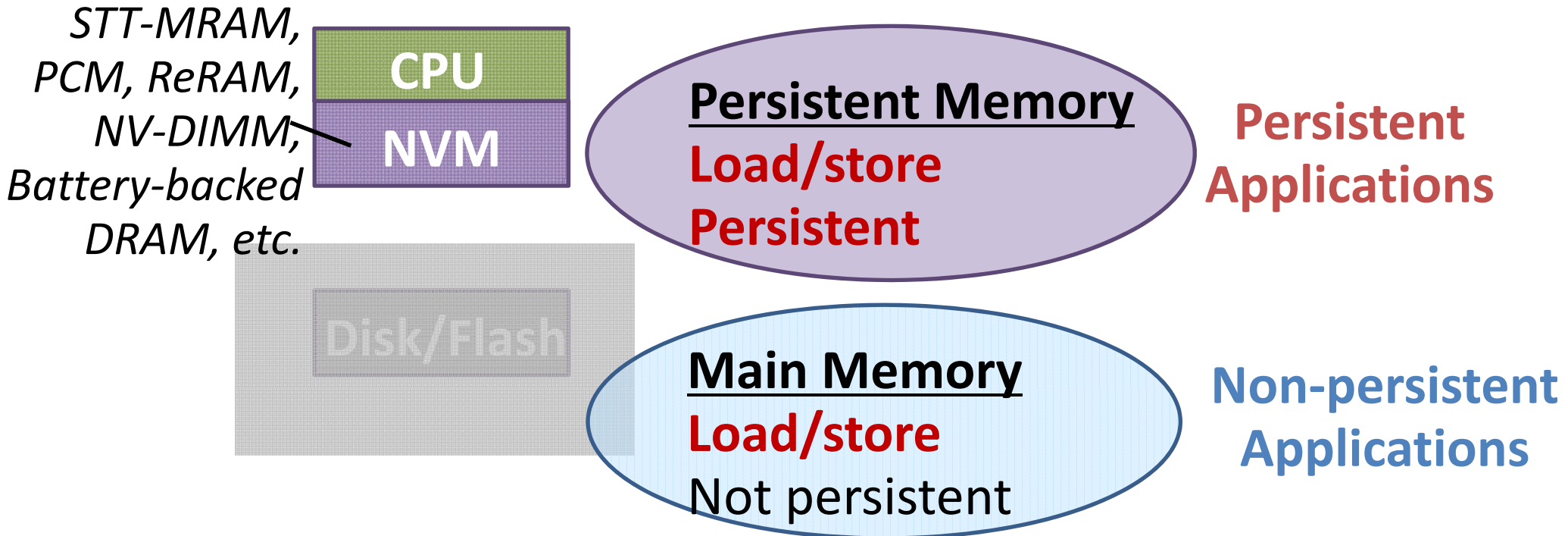
**CPU**

**DRAM**

*Page faults*

**Disk/Flash**

**Main Memory**
**Load/store**
Not persistent

**Storage**
Fopen, fread, fwrite
**Persistent**

# New Design Opportunity with NVMs

*STT-MRAM, PCM, ReRAM, NV-DIMM, Battery-backed DRAM, etc.*

**CPU**

**NVM**

Disk/Flash

**Persistent Memory**
**Load/store**
**Persistent**

**Storage**
Fopen, fread, fwrite
**Persistent**

**Persistent Applications**

**Examples applications**
Databases, file systems, key-value stores
*(In-memory data structures can immediately become permanent)*

# New Design Opportunity with NVMs

*STT-MRAM, PCM, ReRAM, NV-DIMM, Battery-backed DRAM, etc.*

**CPU**

**NVM**

Disk/Flash

**Persistent Memory**
**Load/store**
**Persistent**

**Persistent Applications**

**Main Memory**
**Load/store**
Not persistent

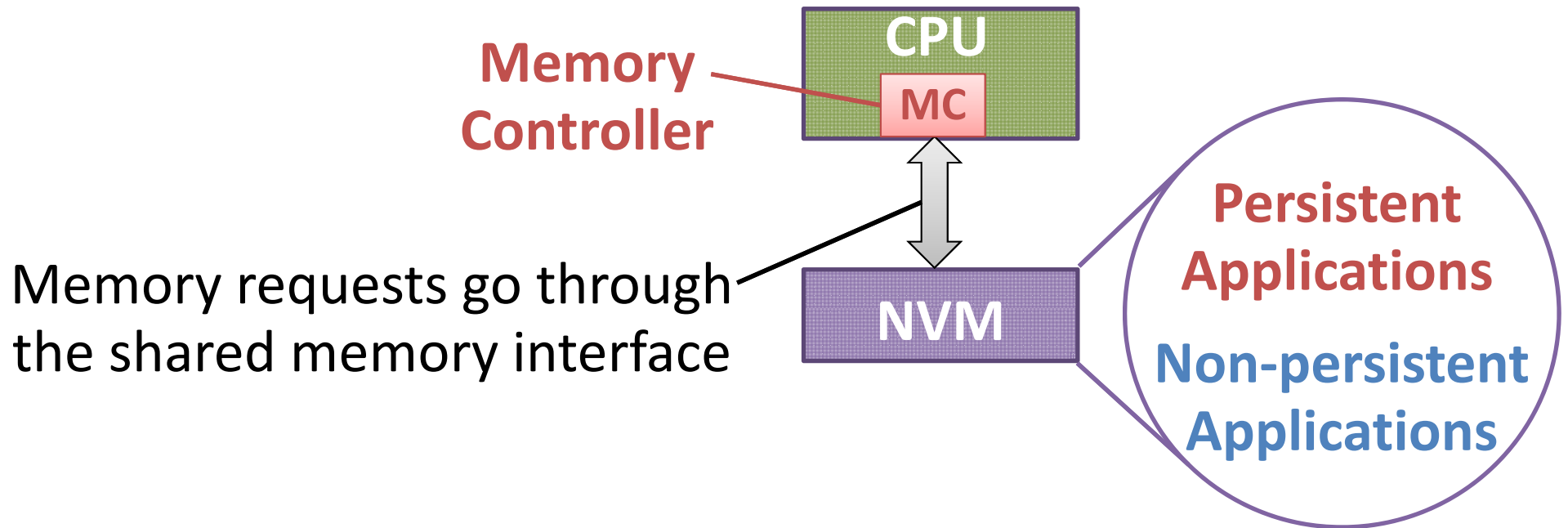**Non-persistent Applications**

**New use case of NVM:**
concurrently running two types of applications
[Kannan + HPCA'14, Liu + ASPLOS'14, Meza + WEED'14]
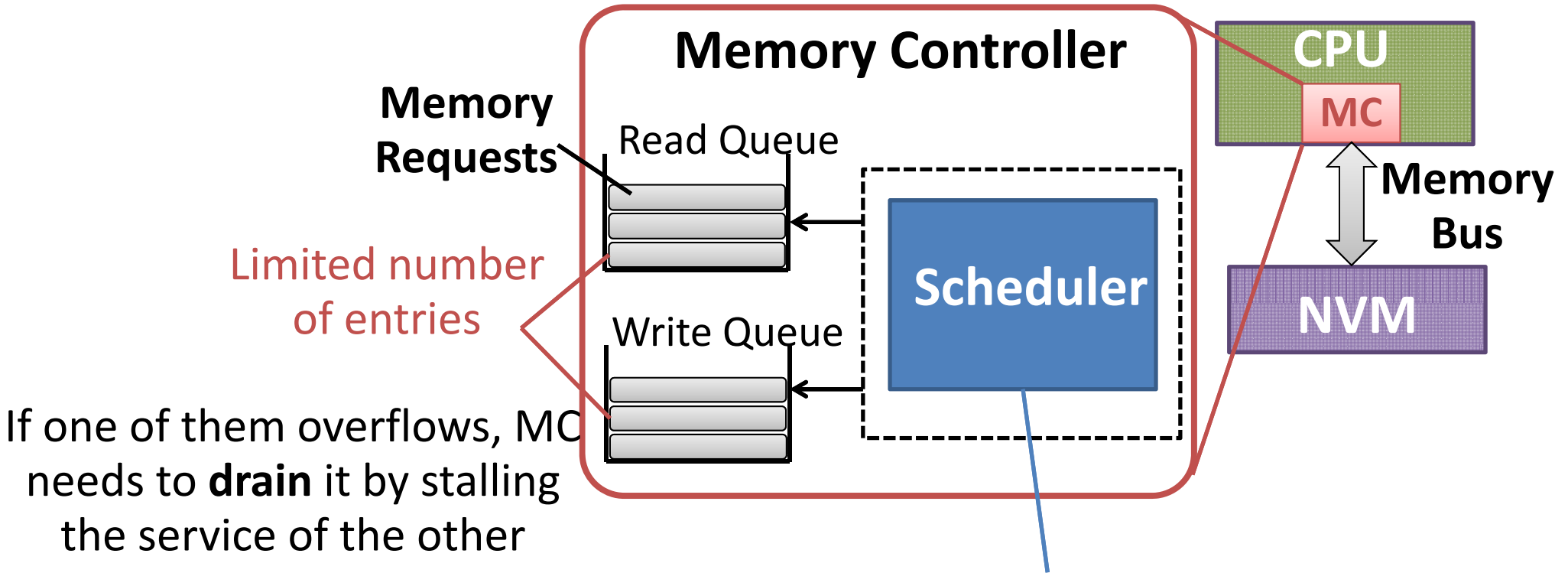
# Focus of Our Work:
# Memory Controller Design

**Memory Controller**

**CPU**

**MC**

Memory requests go through the shared memory interface

**NVM**

**Persistent Applications**

**Non-persistent Applications**

**Fair and High-Performance Memory Control**

# Why Another Memory Control Scheme?

Reads
...
BLISS [ICCD'14]
SMS [ISCA'12]
TCM [MICRO'10]
ATLAS [HPCA'10]
PAR-BS [ISCA'08]
STFM [MICRO'07]
FR-FCFS [ISCA'00]

Writes

**Careful control over writes to guarantee data persistence**

**Persistent Memory**

# Memory Controller



**Memory Requests**

**Memory Controller**

Read Queue

Scheduler

Write Queue

CPU

MC

**Memory Bus**

NVM

Limited number of entries

If one of them overflows, MC needs to **drain** it by stalling the service of the other

Determine which requests can be sent on the memory bus to be serviced

**Why conventional memory control schemes are inefficient in persistent memory systems**

**How to design fair and high-performance memory control in this new scenario**

# Assumptions and Design Choices

## Conventional memory control schemes

| *Assumptions* | *Design choices* |
|---|---|

**1. Reads are on the critical path of application execution**

*(Application execution is read-dependent)*

**1. Prioritize reads over writes**

**2. Applications are usually read-intensive**

**2. Delay writes until they overflow the write queue**

*(Infrequent write queue drains)*

**These assumptions no longer hold in persistent memory, which needs to support data persistence**
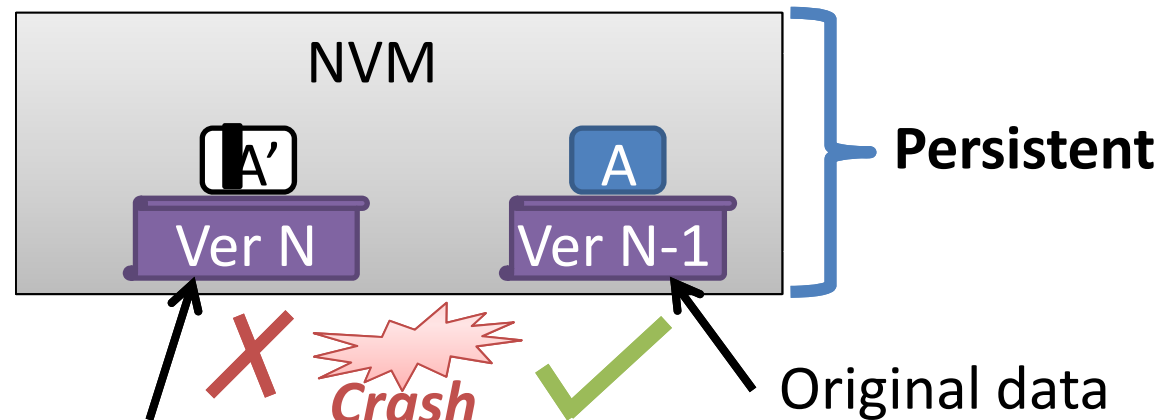
*(Data consistency when the system suddenly crashes or loses power)*

**Mechanisms: multiversioning and write-order control**

# Implication of Multiversioning

Updates of
data structure A
have multiple
write requests

NVM

*Crash*

A

Persistent

# Implication of Multiversioning



**The two versions are not updated at the same time**

**Significantly increasing write traffic –**
Two writes with each one data update

[Volos+ ASPLOS'11, Coburn+ ASPLOS'11, Condit+ SOSP'09, Venkataraman+ FAST'11]

# Assumptions and Design Choices

## *Assumptions*

*Design decisions*

1. Reads are on the critical path of application execution

1. Prioritize reads over writes

2. Applications are usually ~~read-intensive~~

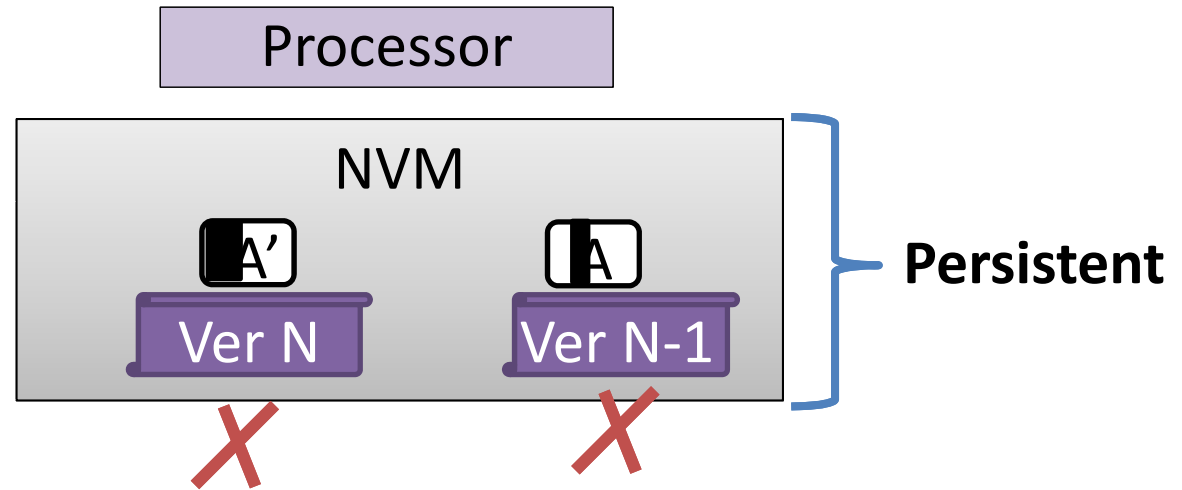**Persistent applications are usually write-intensive**

*Infrequent write queue drains*

# Implication of Write-order Control

*The two versions*
*are not supposed*
*to be updated*
*at the same time,*
*issued in order:*
$A' = \{A'_1, A'_2, A'_3\}$
$A = \{A_1, A_2, A_3\}$
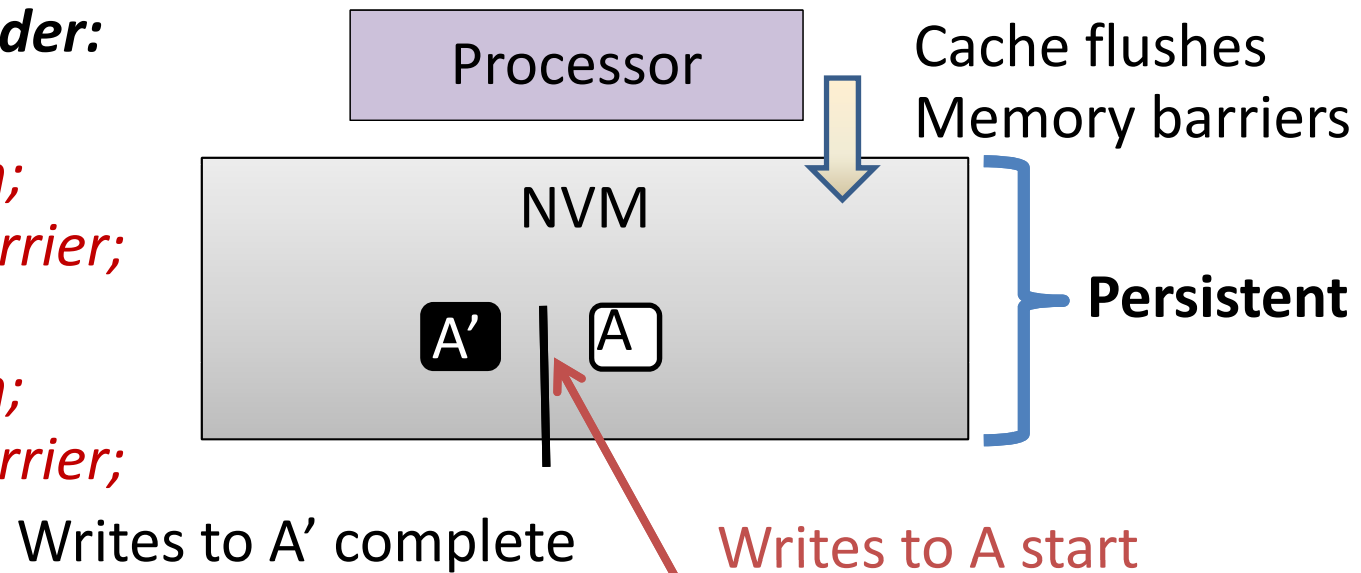
*Reordered by caches and MCs:*
$A_2, A'_2, A'_1, A_1, A_3$

**Crash**

Processor

NVM

A'

A

Ver N

Ver N-1

**Persistent**

# Implication of Write-order Control

*Issued in order:*
*Update A',*
*Cache flush;*
*Memory barrier;*
*Update A ,*
*Cache flush;*
*Memory barrier;*

Processor

NVM

A' | A

Cache flushes
Memory barriers

**Persistent**

Writes to A' complete

Writes to A start

## Restrict the ordering of writes arriving at the memory

**Making application execution write dependent –**
Subsequent writes, reads, and computation can all
depend on a previously issued write

[Volos+ ASPLOS'11, Coburn+ ASPLOS'11, Condit+ SOSP'09, Venkataraman+ FAST'11]

# Assumptions <span style="color:#cccccc">and Design Choices</span>

## *Assumptions*

**1. Reads are on the critical path of application execution**

*(Application execution is read-dependent)*

**Writes are also on the critical execution path**

*(Application execution is write-dependent)*

# Assumptions and Design Choices

*Assumptions*

1. During the ...
application execution

1. Writes are also on the critical execution path

2. Persistent applications are usually write-intensive

**Assumptions**

**Design choices**

1. **Writes are also on the critical execution path**

1. Prioritize rea...

**Unfairness**

2. **Persistent applications are usually write-intensive**

2. Delay writ... they overflow the w... queue

**Performance Degradation**

**Frequent write queue drains**

*Why* ▶

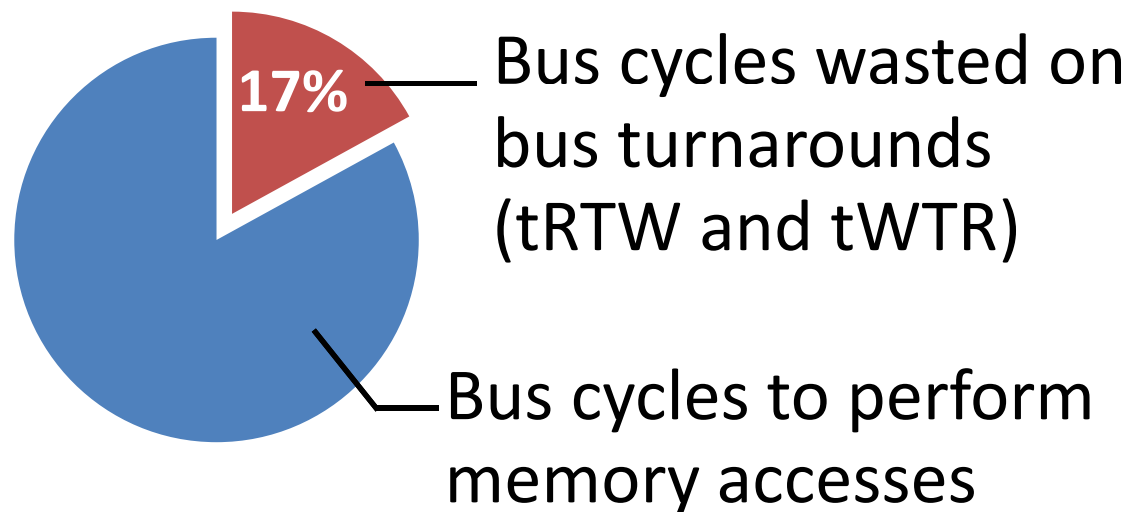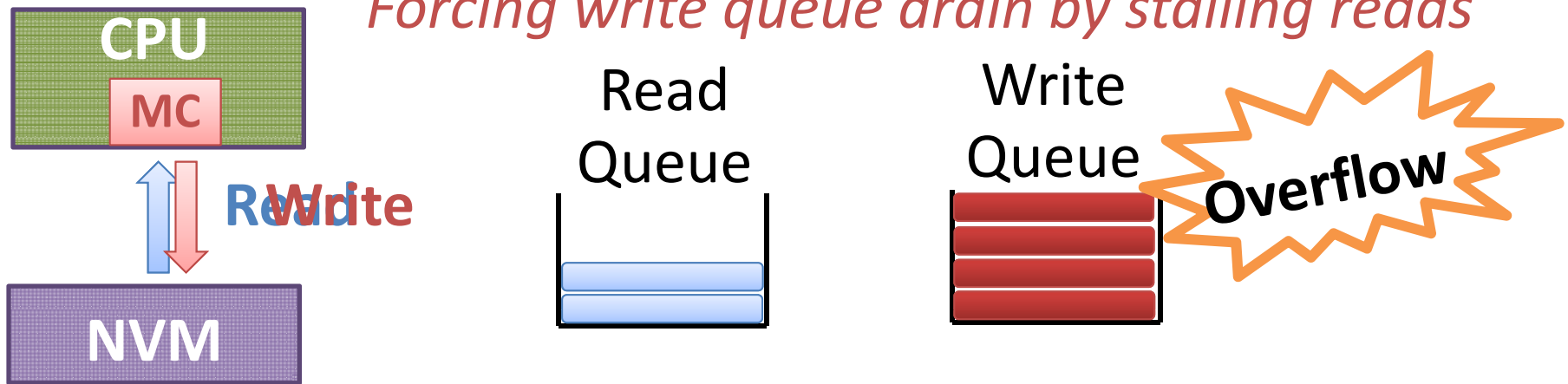*Frequent stall reads to drain the write queue, frequently switch between servicing reads and servicing writes*

# Bus Turnaround Overhead

*tRTW ~ 7.5ns*   *tWTR ~ 15ns*

*[Kim + ISCA'12]*

*Forcing write queue drain by stalling reads*

CPU

MC

NVM

Write Read

Read Queue

Write Queue

Overflow

17% — Bus cycles wasted on bus turnarounds (tRTW and tWTR)

Bus cycles to perform memory accesses

# Assumptions and Design Choices

## *Assumptions*

**1. Writes are also on the critical execution path**

**2. Persistent applications are usually write-intensive**

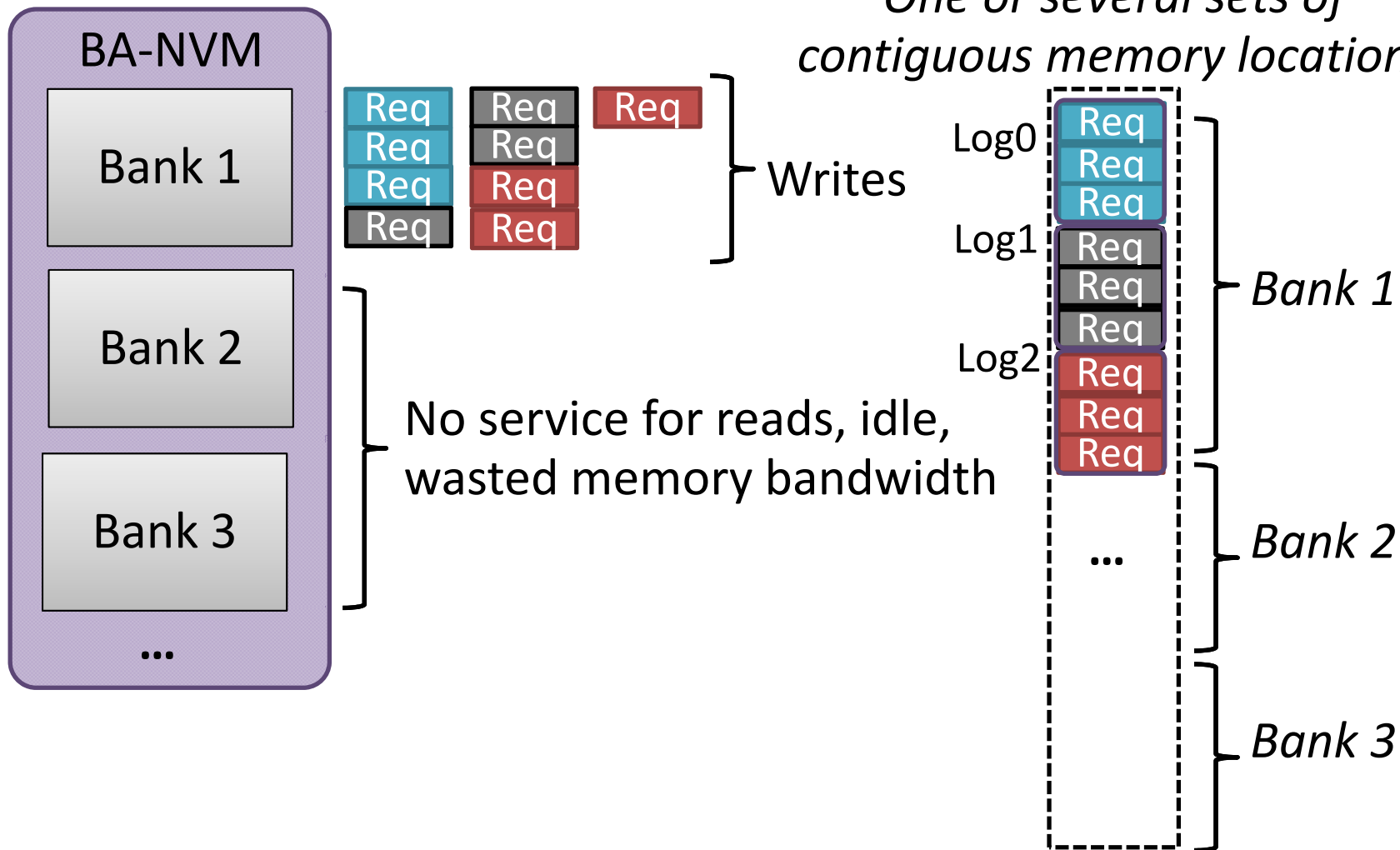**3. Writes in persistent memory have low bank-level parallelism (BLP)**

# Low Bank-level Parallelism (BLP)

*Stalling reads for a long time*
*while the bus is servicing writes to persistent memory*

**A Log  (Spans Multiple Banks)**
*One or several sets of*
*contiguous memory locations*

BA-NVM

Bank 1

Bank 2

Bank 3

...

Req Req Req
Req Req
Req Req
Req Req

Writes

No service for reads, idle,
wasted memory bandwidth

Log0 — Req Req Req
Log1 — Req Req Req
Log2 — Req Req Req

...

Bank 1

Bank 2

Bank 3

# Assumptions and Design Choices

## Assumptions

1. Writes are also on the critical execution path

2. Persistent applications are usually write-intensive

3. Writes of persistent applications have low BLP

## Design choices

? Fairness

? High Performance ?

# Design Principles

**Persistence-Aware Memory Scheduling**
**Minimizing write queue drains and
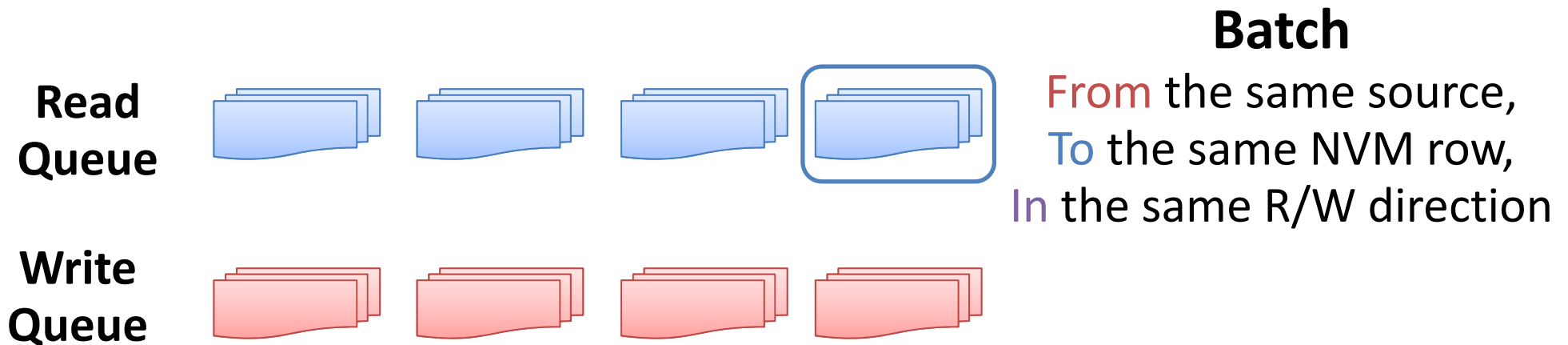bus turnarounds, while ensuring fairness**

**Persistent Write Striding**
**Increasing BLP of writes to persistent memory
to fully utilize memory bandwidth**

# Persistence-aware Memory Scheduling

Minimizing write queue drains and bus turnarounds, while ensuring fairness

*Problem: when to switch between*
*servicing read batches and write batches*
**Low bus turnaround overhead, risk frequent write queue drains**

**Read Queue**

**Write Queue**

**Batch**

From the same source,
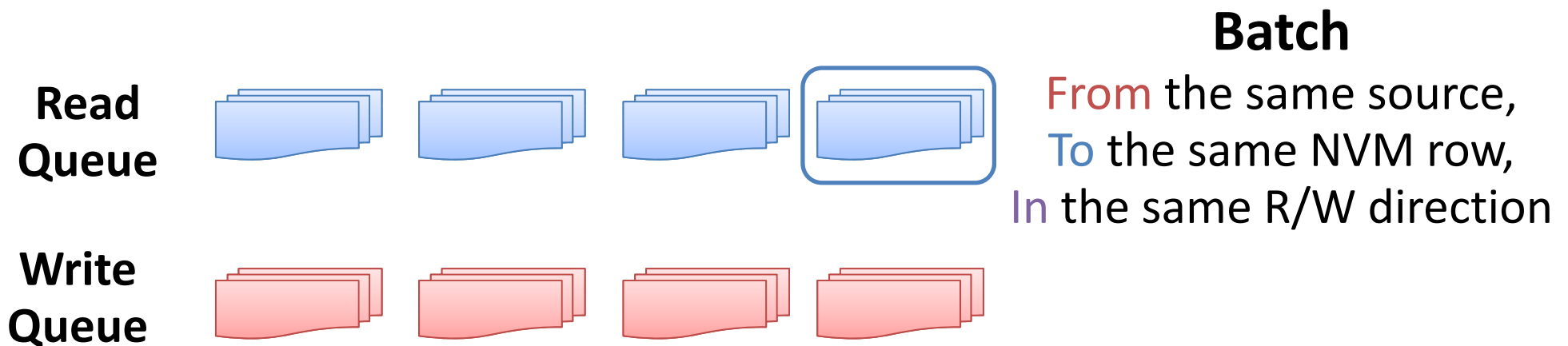To the same NVM row,
In the same R/W direction

# Persistence-aware Memory Scheduling

Minimizing write queue drains and bus turnarounds,
while ensuring fairness

*Problem: when to switch between
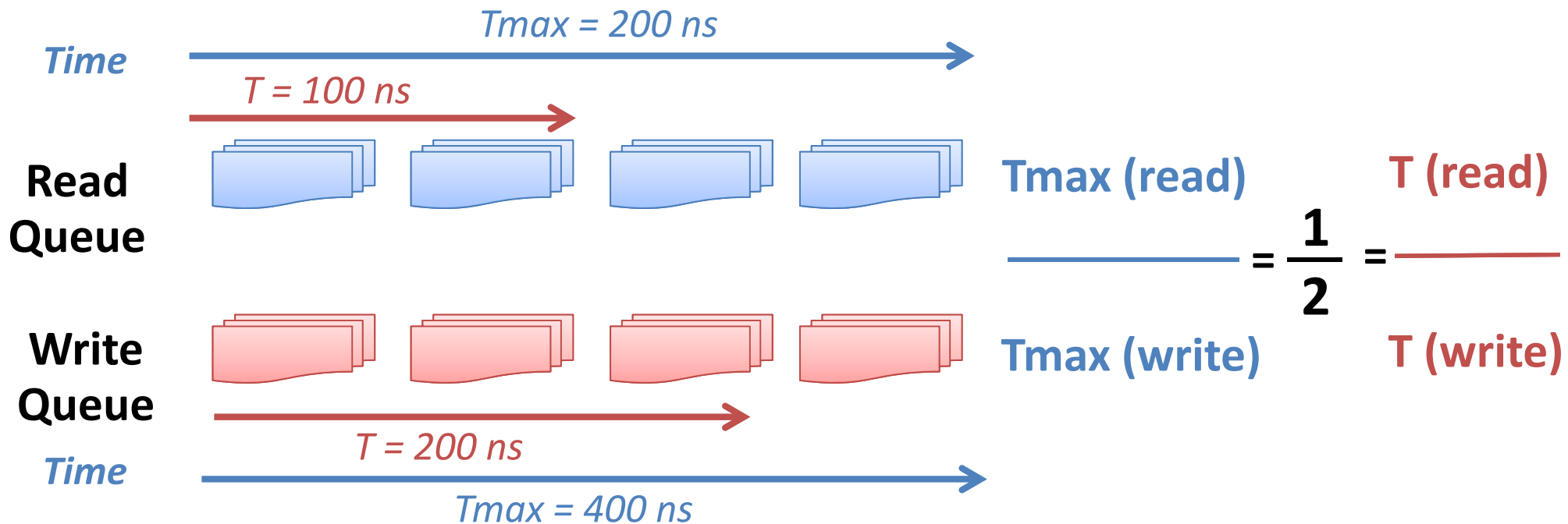servicing read batches and write batches*

**Less likely to starve reads and writes, higher bus turnaround overhead**

**Read Queue**

**Write Queue**

**Batch**

From the same source,
To the same NVM row,
In the same R/W direction

# Persistence-aware Memory Scheduling

Minimizing write queue drains and bus turnarounds, while ensuring fairness

*Key idea 1: balance the amount of time spent in continuously servicing reads and writes*

Time $\longrightarrow$ Tmax = 200 ns

$T = 100$ ns

**Read Queue**

Tmax (read)   T (read)

$$\frac{\text{Tmax (read)}}{\text{Tmax (write)}} = \frac{1}{2} = \frac{\text{T (read)}}{\text{T (write)}}$$

**Write Queue**

Tmax (write)   T (write)

$T = 200$ ns

Time $\longrightarrow$ Tmax = 400 ns

# Persistence-aware Memory Scheduling

Minimizing write queue drains and bus turnarounds, while ensuring fairness

*Key idea 2: Time to service read batches and write batches is JUST long enough*

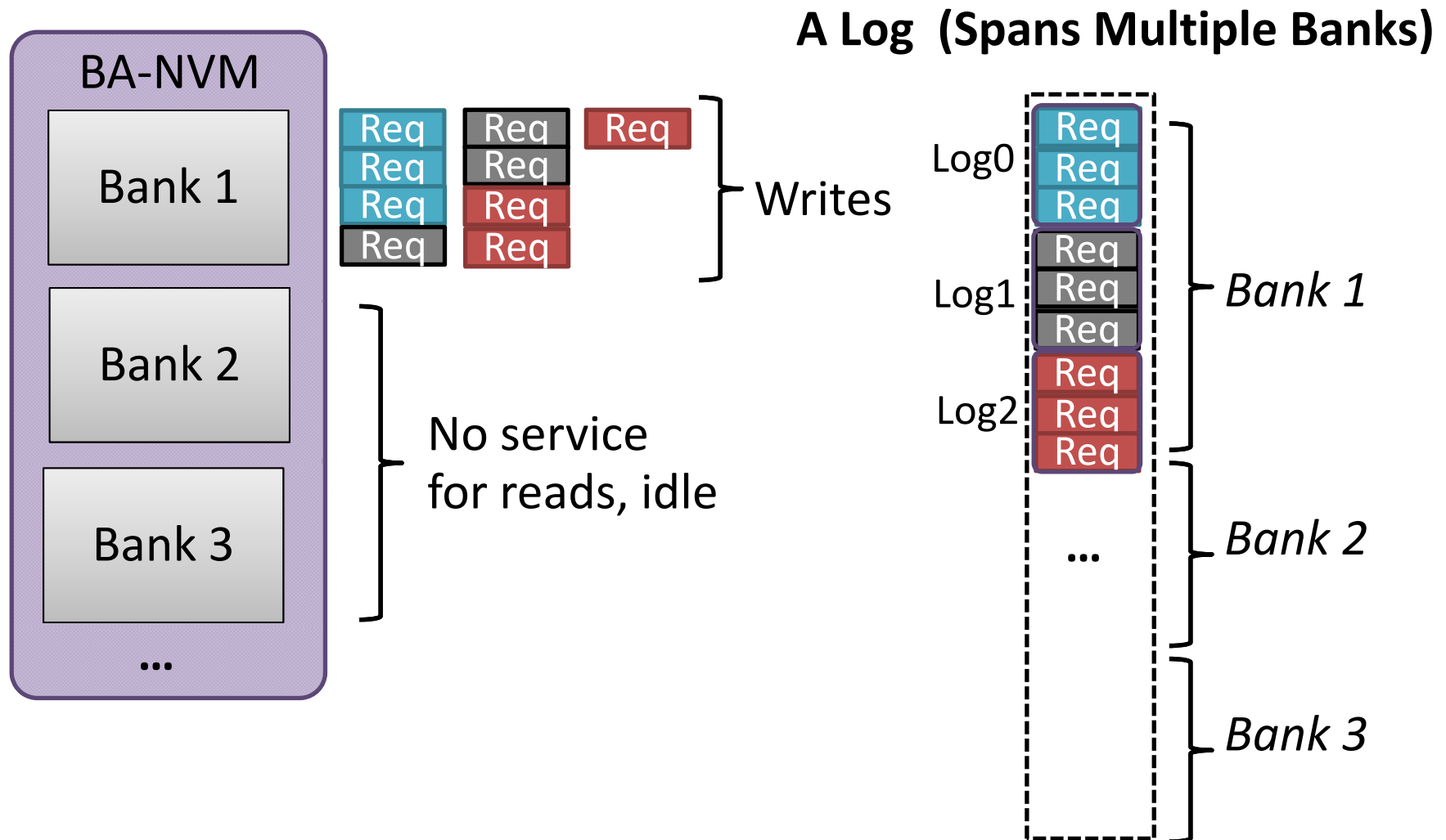$$\frac{tRTW + tWTR}{T(\text{read batches}) + T(\text{write batches})} < \mu$$

*User-defined bus turnaround overhead*

*Pick the choice with the shortest times*

# Persistent Write Striding

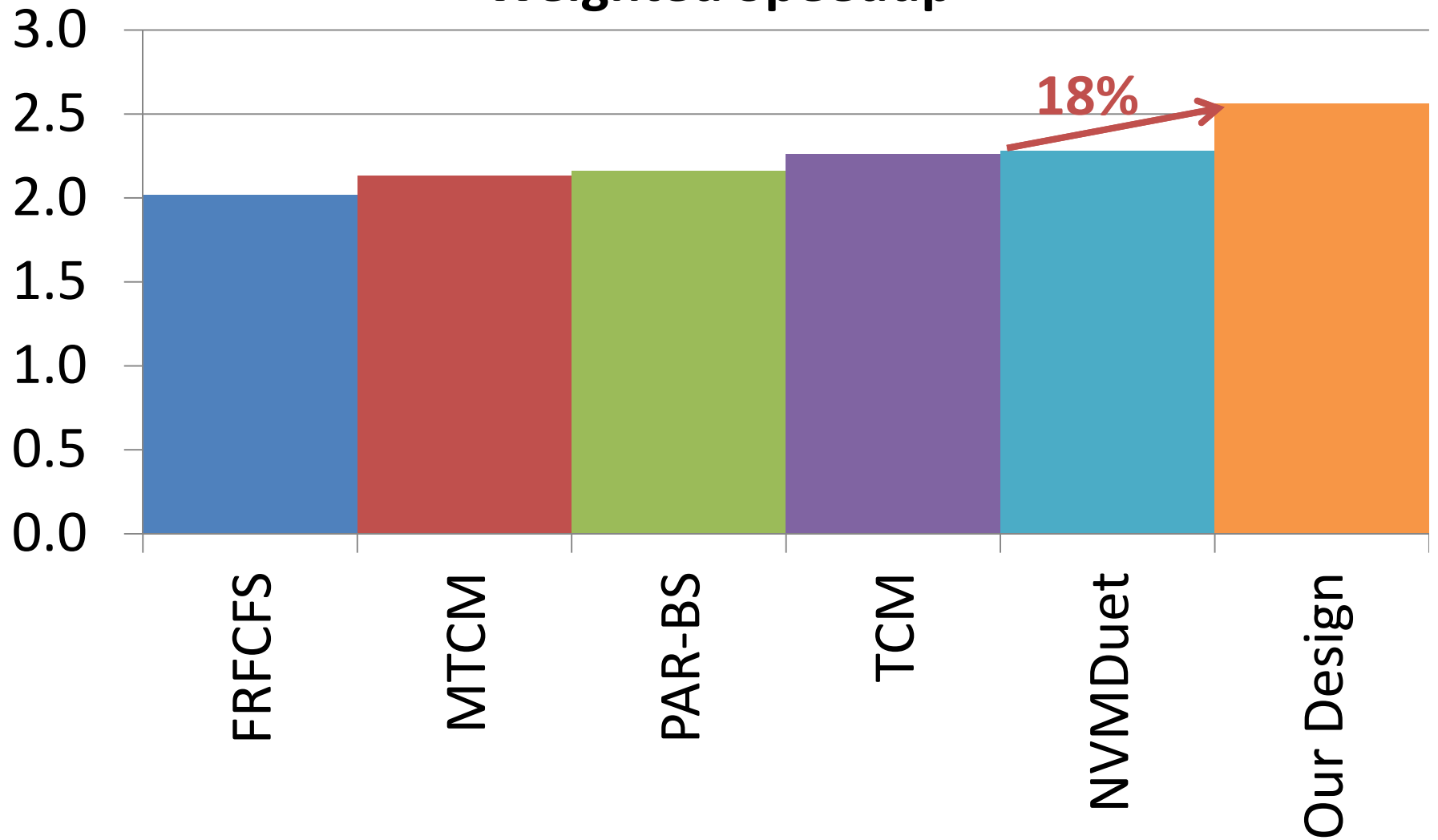Increasing BLP of persistent writes to fully utilize memory bandwidth

## BA-NVM

Bank 1

Bank 2

Bank 3

...

Req Req Req Req } Req Req Req Req } Req Req Req } Writes

No service for reads, idle

## A Log  (Spans Multiple Banks)

Log0 — Req Req Req

Log1 — Req Req Req

Log2 — Req Req Req

} Bank 1

... } Bank 2

} Bank 3

# Persistent Write Striding

Increasing BLP of persistent writes to fully utilize memory bandwidth

*Key idea: stride writes by an offset to remap them to different banks*

**A Log  (Spans Multiple Banks)**

BA-NVM

Req  Req  Req

**Benefit: no bookkeeping required**

Req  Req

Req
Req
Req
Req

Req

Serviced in parallel

Bank 2

Bank 3

...

Log0   Req
       Req

Log K

**Stride + offset**

Log1   Req
       Req
Log2   Req
       Req

Bank 2

**Stride + offset**

Req

Bank 3

...

# Experimental Setup

- ## Simulator
  - McSimA+ [Ahn+, ISPASS'13] (modified)

- ## Configuration
  - Four-core processor, eight threads
  - Private caches: L1/L2, SRAM, 64KB/256KB per core
  - Shared last-level cache: L3, SRAM, 2MB per core
  - Main memory: STT-MRAM DIMM (8GB)

- ## Benchmarks
  - 7 non-persistent applications
  - 3 persistent applications

  18 Workloads
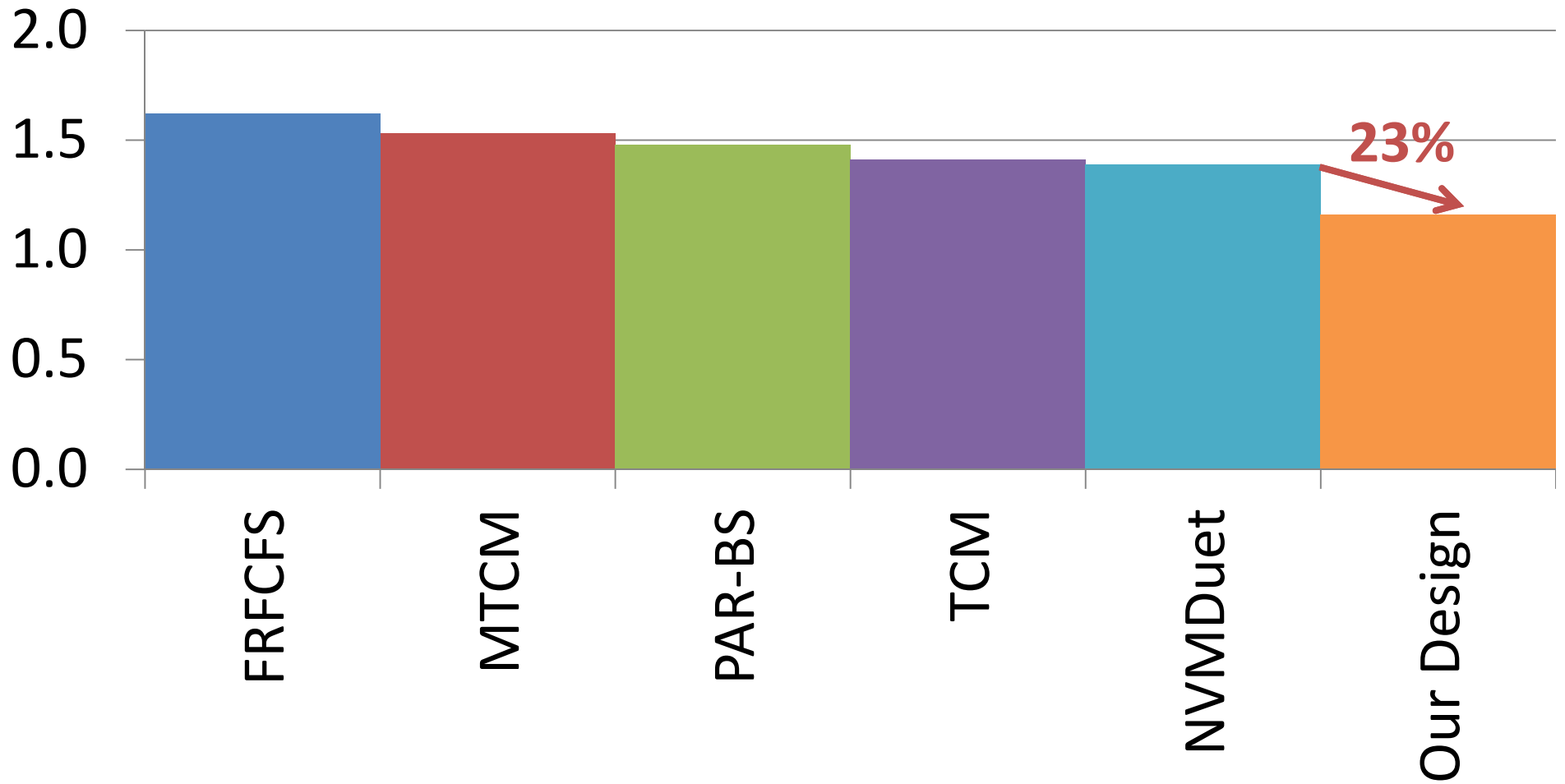
- ## Metrics
  - Weighted speedup & maximum slowdown
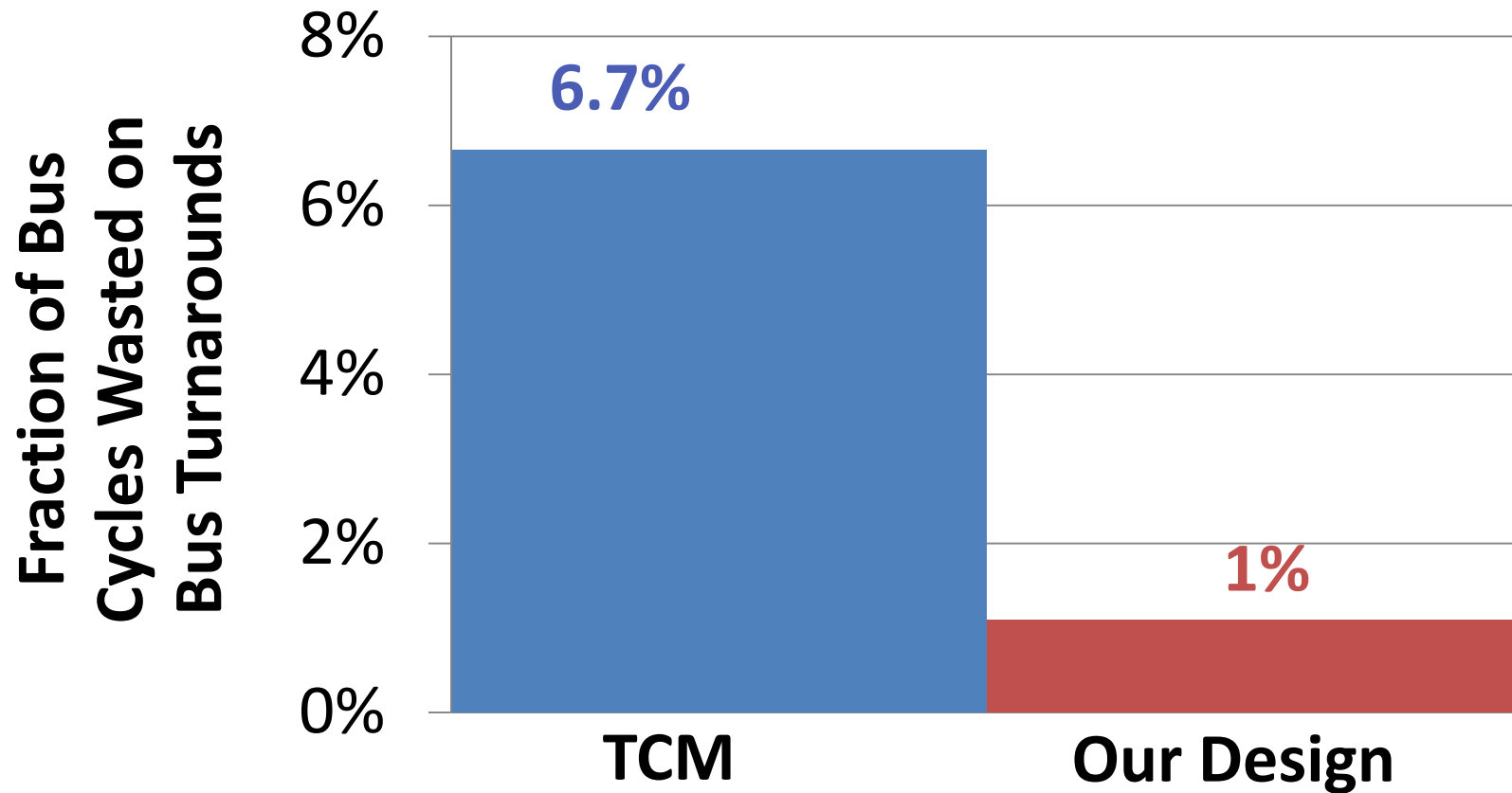
# Performance



Weighted Speedup

# Fairness

# Bus Turnaround Overhead



*(The worst case of previous designs: 17%)*

# Other Results

Sensitivity study on various

- NVM Latencies
- Row-buffer sizes
- Number of threads

# Summary

Reads ⭐ Writes

FIRM [MICRO'14]
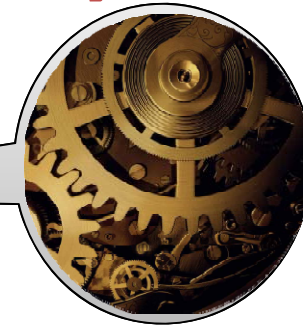
...

BLISS [ICCD'14]

SMS [ISCA'12]

TCM [MICRO'10]

ATLAS [HPCA'10]

PAR-BS [ISCA'08]

STFM [MICRO'07]

FR-FCFS [ISCA'00]

Persistent Memory

**Design principles:**

**Persistence-aware memory scheduling**

**Persistent write striding**

**Fairness**

**High Performance**

# Thank you!

# FIRM: Fair and HIgh-PerfoRmance Memory Control for Persistent Memory Systems

Jishen Zhao       Onur Mutlu       Yuan Xie